



## 1. ❓❓

KB Enterprise/Department/Personal/Agent  
user\_id  
RAG

/

- chunk **100** **1000**
- 70% chunk public**
- kb\_id
- embedding **768**
- chunk **600 tokens**
- Hybrid **95% BM25 + kNN**
- rerank LLM rerank**

## 2. ❓❓

- BM25 + kNN/HNSW + Hybrid + Rerank**
- /
- user\_id 0**
- 100 chunk **1000 chunk**

## 3. ❓❓❓

### 3.1 ❓❓❓

- scope\_id**
- MVP** **100 chunk** **scope\_id**
- (Scale)**

- $\diamond$  public  $\diamond$  private  $\diamond\diamond\diamond\diamond\diamond\diamond\diamond\diamond\diamond$
  - public  $\diamond\diamond\diamond\diamond$  ACL filter
  - private  $\diamond\diamond\diamond\diamond$  ACL pre-filter  $\diamond\diamond\diamond=0\diamond$

4. Hybrid 95% + rerank 8888888888chunk 600 tokens8888888888

## ◆◆ rerank ◆◆◆◆◆

5. 问题/OCR/ASR 文本嵌入

6.    &

- doc\_type/long/short
  - chunk/keyframe/chunk-level Q rerank
  - doc\_meta / chunk\_meta / / /

# 4.

## 4.1

- `user_id`
  - `chunk` `scope_id`

## 4.2 scope\_id ◊◊

- `public_all` 70% chunk
  - `/agent scope id`

## 4.3 ACL

## 4.4 🔑=0 🔑

- private ES terms scope\_id in user\_scopes filter context
- scope

## 5. 🔑?

### 5.1 🔑?

- kb\_chunks\_v1
- public/private scope\_id
- terms scope\_id in user\_scopes

?

public ACL filter Scale

### 5.2 public/private 🔑?

- kb\_chunks\_public\_v1 scope\_id=public\_all
- kb\_chunks\_private\_v1 scope\_id!=public\_all

?

- public BM25 + public kNN ACL filter
- private BM25 + private kNN ACL filter  
+ rerank

## 6. 🔑? Ingestion

### 6.1 🔑?

- PDF/Office/HTML/Markdown/IM
- ?
- chunk 600 tokens overlap 80 120 400 800 tokens
- embedding 768 embedding\_model/embedding\_version

## 6.2

- ②②②OCR → ②② → chunk → embedding
  - ②②②ASR → transcript → chunk②②②+token②②②→ embedding
  - ②②②②②/ASR/frame → transcript → chunk②②②/②②②②→ embedding
  - ②②②②②②② start\_ms/end\_ms

## 6.3

- content\_hash SHA1/xxhash
  - quality\_score
  - chunk\_id ES \_id

## 7. Mapping JSON

```
PUT kb_chunks_v1
{
  "settings": {
    "number_of_shards": 3,
    "number_of_replicas": 1,
    "refresh_interval": "10s",
    "analysis": {
      "analyzer": {
        "zh_smart": {
          "type": "ik_smart"
        },
        "zh_max": {
          "type": "ik_max_word"
        },
        "en_std": {
          "type": "standard"
        },
        "pinyin_analyzer": {
          "tokenizer": "standard",
          "filter": [
            "lowercase",
            "my_pinyin"
          ]
        }
      },
      "filter": {
        "my_pinyin": {
          "type": "pinyin",
          "keep_full_pinyin": true,
          "keep_first_letter": true,
          "keep_separate_first_letter": false,
          "keep_joined_full_pinyin": true,
          "keep_original": true,
          "limit_first_letter_length": 16,
          "remove_duplicated_term": true
        }
      }
    }
  }
},
```

```
"mappings": {
    "dynamic": "strict",
    "properties": {
        "kb_id": { "type": "keyword" },
        "kb_type": { "type": "keyword" },

        "scope_id": { "type": "keyword" },

        "doc_id": { "type": "keyword" },
        "chunk_id": { "type": "keyword" },
        "chunk_index": { "type": "integer" },

        "language": { "type": "keyword" },

        "title": {
            "properties": {
                "zh": {
                    "type": "text",
                    "analyzer": "zh_smart",
                    "fields": {
                        "max": { "type": "text", "analyzer": "zh_max" },
                        "pinyin": { "type": "text", "analyzer": "pinyin_analyzer" },
                        "keyword": { "type": "keyword", "ignore_above": 256 }
                    }
                },
                "en": { "type": "text", "analyzer": "en_std" }
            }
        },
        "content": {
            "properties": {
                "zh": {
                    "type": "text",
                    "analyzer": "zh_smart",
                    "fields": {
                        "max": { "type": "text", "analyzer": "zh_max" },
                        "pinyin": { "type": "text", "analyzer": "pinyin_analyzer" }
                    }
                },
                "en": { "type": "text", "analyzer": "en_std" }
            }
        }
    }
}
```

```

        },
    },

    "source_type": { "type": "keyword" },
    "doc_type": { "type": "keyword" },

    "created_at": { "type": "date" },
    "updated_at": { "type": "date" },

    "quality_score": { "type": "float" },
    "tags": { "type": "keyword" },
    "content_hash": { "type": "keyword" },

    "start_ms": { "type": "long" },
    "end_ms": { "type": "long" },

    "embedding_model": { "type": "keyword" },
    "embedding_version": { "type": "keyword" },

    "embedding": {
        "type": "dense_vector",
        "dims": 768,
        "index": true,
        "similarity": "cosine"
    }
}
}
}

```

## 8. ◊◊/◊◊/◊◊/◊◊◊◊◊◊

### 8.1 MVP

- shards ◊3~6◊◊◊ 3◊◊◊◊◊◊◊◊
- replicas ◊1◊◊◊◊◊◊◊◊◊◊ 0◊
- refresh\_interval ◊10s

## 8.2 Scale

- public ◇ 700 ◇ ◇ ◇ shards=12 ◇ 8~16 ◇ ◇ replicas=1
- private ◇ 300 ◇ ◇ ◇ shards=6 ◇ 4~8 ◇ ◇ replicas=1

## 8.3 bulk ◇ ◇

- ◇ ◇ 5~15MB ◇ ◇ 1000~5000 docs/◇ ◇ ◇ doc ◇ ◇ ◇ ◇
- ◇ ◇ 2~8 ◇ ◇ ◇ ◇ ◇ ◇ ◇ ◇ ◇
- ◇ ◇ ◇ ◇ refresh=-1 → ◇ ◇ ◇ ◇ ◇ ◇ → ◇ ◇ ◇ ◇ ◇ segment

## 9. ◇ ◇ ◇ ◇

### 9.1 ◇ ◇ ◇ ◇

- BM25 size ◇ 200 ◇ 150~400 ◇
- kNN ◇ k=150 ◇ 100 250 ◇ ◇ num\_candidates=2000 ◇ 1000 6000 ◇
- ◇ ◇ topM ◇ 150~200
- rerank topR ◇ 60~120 ◇ ◇ ◇ 100 ◇
- ◇ ◇ topK ◇ 10/20/40 ◇ ◇ ◇ 20 ◇

## 9.2 ◊◊◊ BM25 DSL◊◊◊◊◊◊◊

```
POST kb_chunks_v1/_search
{
  "size": 200,
  "track_total_hits": false,
  "timeout": "250ms",
  "query": {
    "bool": {
      "filter": [
        { "terms": { "scope_id": ["public_all","dept_finance","project_abc"] } }
      ],
      "must": [
        {
          "multi_match": {
            "query": "搜索关键词",
            "fields": ["title.zh^3","content.zh","title.en^2","content.en"],
            "type": "best_fields"
          }
        }
      ]
    }
  }
}
```

## 9.3 ◊◊◊ kNN DSL◊◊◊ pre-filter◊◊◊◊◊

```
POST kb_chunks_v1/_search
{
  "size": 150,
  "track_total_hits": false,
  "timeout": "400ms",
  "knn": {
    "field": "embedding",
    "query_vector": [0.01, 0.02, "..."],
    "k": 150,
    "num_candidates": 2000,
    "filter": {
      "bool": {
        "filter": [
          { "terms": { "scope_id": ["public_all", "dept_finance", "project_abc"] } }
        ]
      }
    }
  }
}
```

## 9.4 Hybrid ◊◊◊RRF◊

- ◊ BM25 ◊ kNN ◊◊◊◊◊◊◊◊◊◊◊ chunk\_id◊
- RRF◊ score =  $\sum 1/(k_0 + \text{rank}_i)$  ◊◊ k<sub>0</sub>=60
- ◊◊◊◊◊ topM=200 → ◊◊◊ rerank topR=100

## 10. ◊◊◊◊◊

### 10.1 topK ◊◊◊◊◊

- topK=10◊◊◊◊◊/◊◊◊◊◊
- topK=20◊◊◊◊◊/◊◊◊◊◊
- topK=40◊◊◊◊◊/◊◊◊◊◊/◊◊◊◊◊

# 10.2



# 10.3 chunk

- topK doc chunkchunk\_index LLM
  - token 8k/16k

## 11. Rerank

- BM25+kNN
  - **LLM rerank** 20~50
  - freshness updated\_at quality\_score  
asr/ocr\_confidence /
  - rerank 200~800ms / topR  
LLM rerank

# 12.    ?    ?    ?

# 12.1

1. ☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐ PDF/Office/HTML/OCR/ASR ☐☐
  2. ☐☐☐/☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐☐

3. Chunk
4. hash
5. Embedding
6. scope\_id
7. public → kb\_chunks\_public\_write → non-public → kb\_chunks\_private\_write

## 12.2

sequenceDiagram participant U as 1. User participant S as 2. SearchService participant ACL as 3. ACL/ScopeCache participant ES as 4. Elasticsearch participant RR as 5. RerankService participant L as 6. LLM/RAG U->>S: 1) query + user\_id (+topK) S->>ACL: 2) GetUserScopes(user\_id) ACL-->>S: 3) scopes (public\_all) par BM25 S->>ES: 4a) BM25 search (filter scopes) ES-->>S: bm25\_hits and kNN S->>ES: 4b) kNN search (filter scopes) ES-->>S: knn\_hits end S->>S: 5) + RRF (topM) S->>RR: 6) rerank(topR passages/chunks) RR-->>S: rerank\_scores S->>S: 7) + chunk S->>L: 8) topK chunks/ passages + L-->>S: 9) answer + citations S-->>U: 10) answer + citations

1. query user\_id topK
2. scopes public\_all
3. BM25 kNN scopes pre-filter
4. + RRF topM
5. topR rerank cross-encoder LLM rerank
6. chunk
7. topK LLM/RAG

### public/private

sequenceDiagram participant U as 1. User participant S as 2. SearchService participant ACL as 3. ACL/ScopeCache participant ESP as 4a/4b. ES Public participant ESR as 4c/4d. ES Private participant RR as 6. RerankService participant L as 8. LLM/RAG U->>S: 1) query + user\_id (+topK) S->>ACL: 2) GetUserScopes(user\_id) ACL-->>S: 3) scopes (public\_all) par Public S->>ESP: 4a) Public BM25 (no ACL) ESP-->>S: bm25\_pub S->>ESR: 4b) Public kNN (no ACL) ESP-->>S: knn\_pub and Private S->>ESR: 4c) Private BM25 (filter scopes) ESR-->>S:

>>S: bm25\_pri S->>ESR: 4d) Private kNN (filter scopes) ESR-->S: knn\_pri end S->>S: 5) + RRF (topM) S->>RR: 6) rerank(topR) RR-->S: rerank\_scores S->>S: 7) + (chunk) S->>L: 8) topK + L-->S: 9) answer + citations S-->>U: 10) answer + citations  
 ? ? ? ? ? ? ?

1. query + user\_id topK
2. scopes public\_all
3. BM25 ACL kNN ACL
4. BM25 ACL filter 4d kNN ACL filter
5. chunk\_id RRF topM
6. topR rerank cross-encoder LLM rerank
7. chunk
8. topK LLM/RAG
9. answer + citations

Scale ES public/private public ACL filter private  
 filter

## 14. 🔎

- OCR/ASR
- - /quality\_score update content\_hash
  - /chunk doc\_id → embedding →
  - /chunk\_id
- - chunk\_id \_id content\_hash
  - bulk refresh\_interval 10s refresh
  - doc\_id chunk
-

- `update nDCG/Recall`
  - `kb_chunks_backup version`

ANSWER

sequenceDiagram participant U as 1. Producer participant V as 2. Parser/Chunker participant S as 3. IngestService participant ES as 4. Elasticsearch participant OB as 5. Observability U->>V: 1) `doc_id` V->>V: 2) `embedding` embedding V-->>S: 3) chunk + meta + content\_hash S->>ES: 4) `delete_by_query(doc_id)` S->>ES: 5) bulk index/update (`chunk_id=_id`) ES-->>S: 6) ack S->>ES: 7) optional refresh S-->>OB: 8) `update` /`chunk_id=_id`



# 15.

?

A horizontal row of twelve identical diamond-shaped frames, each with a black border and a white center. Inside each frame is a black question mark. The frames are evenly spaced and extend across most of the width of the page.

1. Operator `/`
  2. IngestService `update` `is_deleted=true, deleted_at`  
`delete_by_query``doc_id/chunk_id`
  3. `scope_id` `ACL`
  4. `tombstone`  
`scroll_size`
  5. `doc_id`
  6. `segment`

# 17.

- Mapping embedding / ILM
  - +

- alias alias
- kb\_chunks\_public\_write\_vX / kb\_chunks\_private\_write\_vX  
kb\_chunks\_public / kb\_chunks\_private

sequenceDiagram participant OP as 1. Operator/ControlPlane participant IG as 2. IngestService participant NEW as 3. NewIndex (\_write vX) participant OLD as 4. OldIndex (oldIndex) participant AL as 5. AliasSwitch participant OB as 6.

Observability OP->>IG: 1) mapping/settings  
IG->>NEW: 2) mapping/settings  
IG->>NEW: 3) bulk refresh=-1, replicas=0 opt  
IG->>OLD: 4a) end  
IG->>NEW: 4b) end  
IG->>IG: 5) /  
IG->>AL: 6) /  
AL-->>IG: 7) ack  
IG->>NEW: 8) refresh\_interval/replicas  
IG-->>OB: 9) /  
IG->>OLD: 10) /

refresh\_interval/replicas

refresh\_interval/replicas

- shards/replicas refresh\_interval -1
- Mapping/Settings \*\_write vX analyzer/
- replica refresh=-1 5~15MB bulk
- IngestService
- doc\_id/chunk scope
  - BM25/kNN/Hybrid = 0
- ingest refresh\_interval replica=1
- 24~72 force merge segment

reject GC nDCG/Recall

- reject GC nDCG/Recall
- /

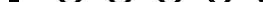
# 18.



A horizontal row of seven empty diamond-shaped boxes, each containing a question mark, used for user input.

sequenceDiagram participant OP as 1. Operator participant IG as 2. IngestService participant P as 3. Parser/Chunker participant ES as 4. Elasticsearch participant OB as 5. Observability OP->>IG: 1) ?/?/?/?/?/?/? (doc\_id scope/?/?/? ) IG->>P: 2) ?/?/?/?/?/?/?/? P-->>IG: 3) ?? chunk+meta+embedding IG->>ES: 4) delete\_by\_query(doc\_id/chunk\_id) ?? IG->>ES: 5) bulk ?? refresh? IG->>ES: 6) scope/TTL/?/?/?/?/?/?/? IG-->>OB: 7) ??/?/?/?/?/?/?/? ?/?/?/?/?

1. doc\_id scan updated\_at
  2. Parser/Chunker chunk meta content\_hash embedding/
  3. doc\_id delete\_by\_query chunk\_id scope
  4. Bulk refresh
  5. ACL
  6. /ACL/ scope /ACL/ user\_scopes
  - doc\_id + scope

**19.**    

## 19.1

- BM25/kNN ◊◊◊ P50/P95/P99◊

- ⓘ ⓘ ⓘ ⓘ ⓘ rerank ⓘ ⓘ
  - ⓘ ⓘ ⓘ ⓘ ⓘ Recall@K ⓘ nDCG ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ ⓘ = 0 ⓘ

# 19.2

- query >= 1~2
  - 20 QPS 100 QPS 150 QPS
  - topK=10/20/40 rerank on/off

## 19.3

## ❖ ES reject/❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖ ❖

1. ◊ num\_candidates
  2. ◊ k
  3. ◊ topR ◊ rerank ◊ ◊ ◊ ◊ ◊ ◊
  4. ◊◊◊ BM25-only ◊ ◊ ◊ ◊ ◊ ◊



# 21.     ?     ?     ?



?	?
Cumulative Gain)	
num_candidates	2000 1000~6000 ?
quality_score	?
refresh_interval	MVP 10s ?
replicas	1 0 ?
rerank	cross-encoder LLM ?
RRF (Reciprocal Rank Fusion)	?
scope	scope_id public_all ACL ?
scope_id	?
Shadow index	?
shards	MVP 3~6 3 ?
Tombstone	is_deleted=true ?
topK	RAG 10/20/40 20 ?
topM	BM25+kNN 150~200 rerank ?
topR	rerank 60~120 100 ?
TTL (Time To Live)	scope 1~5 ?
scope_id=0	scope_id ?